



**Europäisches  
Patentamt**

**European  
Patent Office**

**Office européen  
des brevets**

10978 U.S. PTO  
09/992960  
11/14/01

**Bescheinigung**

**Certificate**

**Attestation**

Die angehefteten Unterla-  
gen stimmen mit der  
ursprünglich eingereichten  
Fassung der auf dem näch-  
sten Blatt bezeichneten  
europäischen Patentanmel-  
dung überein.

The attached documents  
are exact copies of the  
European patent application  
described on the following  
page, as originally filed.

Les documents fixés à  
cette attestation sont  
conformes à la version  
initialement déposée de  
la demande de brevet  
européen spécifiée à la  
page suivante.

**Patentanmeldung Nr.    Patent application No.    Demande de brevet n°**

00125014.1

Der Präsident des Europäischen Patentamts;  
Im Auftrag

For the President of the European Patent Office

Le Président de l'Office européen des brevets  
p.o.

**I.L.C. HATTEN-HECKMAN**

DEN HAAG, DEN  
THE HAGUE,    22/10/01  
LA HAYE, LE

**THIS PAGE BLANK (USPTO)**



Europäisches  
Patentamt

European  
Patent Office

Office européen  
des brevets

10978 U.S. PTO  
09/992960  
11/14/01

**Blatt 2 der Bescheinigung  
Sheet 2 of the certificate  
Page 2 de l'attestation**

Anmeldung Nr.:  
Application no.:  
Demande n°: 00125014.1

Anmeldetag:  
Date of filing:  
Date de dépôt: 16/11/00

Anmelder:  
Applicant(s):  
Demandeur(s):  
Sony International (Europe) GmbH  
10785 Berlin  
GERMANY

Bezeichnung der Erfindung:  
Title of the invention:  
Titre de l'invention:  
Method for recognizing speech

In Anspruch genommene Priorität(en) / Priority(ies) claimed / Priorité(s) revendiquée(s)

Staat:  
State:  
Pays:

Tag:  
Date:  
Date:

Aktenzeichen:  
File no.  
Numéro de dépôt:

Internationale Patentklassifikation:  
International Patent classification:  
Classification internationale des brevets:

G10L15/08

Am Anmeldetag benannte Vertragsstaaten:  
Contracting states designated at date of filing: AT/BE/CH/CY/DE/DK/ES/FI/FR/GB/GR/IE/IT/LI/LU/MC/NL/PT/SE/TR  
Etats contractants désignés lors du dépôt:

Bemerkungen:  
Remarks:  
Remarques:

**THIS PAGE BLANK (USPTO)**

MÜLLER & HOFFMANN - PATENTANWÄLTE

European Patent Attorneys - European Trademark Attorneys

Innere Wiener Strasse 17  
D-81667 München

EPO - Munich  
17

16. Nov. 2000

Attorney's File: 52.744

Gö/Sei/le

Applicant's Case: S00P5131EP00  
PAE00-068HNCE

16.11.2000

Sony International (Europe) GmbH

Kemperplatz 1

D - 10785 Berlin  
GERMANY

---

METHOD FOR RECOGNIZING SPEECH

---

## Description

EPO - Munich  
17

16. Nov. 2000

- 1 The present invention relates to a method for recognizing speech according to claim 1, and in particular to a method for recognizing speech using confidence measures in a process of large vocabulary continuous speech recognition (LVCSR).

5

- In many conventional devices and methods for recognizing speech after recognition of a received utterance or speech phrase an estimation is given on the reliability of the recognized utterance or speech phrase, in particular to enable a decision on whether or not the utterance or speech phrase in question  
10 and its recognized form can be accepted for further processing or has to be rejected and to be exchanged by an utterance or speech phrase to be entered newly by the speaker or user.

- A major drawback of prior art methods for recognizing speech is that the total  
15 computational burden is distributed over the entire received utterance to ensure a detailed and thorough analysis. Therefore, many methods cannot be implemented in small systems or devices, for example in hand-held appliances or the like, as these small systems possess a performance rate which is not sufficient to recognize continuous speech and estimate the reliability of the  
20 recognized phrases when the entire received utterance has to be thoroughly analyzed.

- It is therefore an object of the present invention to provide a method for recognizing speech, in particular in the field of large vocabulary continuous  
25 speech recognition, which can easily be implemented in small dialogue systems and which also gives a robust and reliable estimation on the recognition quality.

- The object is achieved by a method for recognizing speech with the  
30 characterizing features of claim 1. Preferred embodiments of the inventive method for recognizing speech are within the scope of the dependent claims.

- In the method for recognizing speech according to the invention a received utterance is subjected to a recognizing process in its entirety. Further, an only  
35 rough estimation is made on whether or not said received and recognized

- 1 utterance is accepted or rejected in its entirety. Additionally, in the case of  
accepting said utterance it is thoroughly reanalyzed to extract its meaning  
and/or intention. Additionally, based on the reanalysis and its result key-  
phrases and/or keywords are extracted from the utterance essentially being  
5 representative for its meaning.

In contrast to prior art methods for recognizing speech after recognizing the  
utterance in its entirety within a recognizing process an only rough estimate is  
performed describing the reliability of the recognized utterance for necessary  
10 speech phrases. Therefore, only a small burden of estimation and calculation is  
to be focussed on the entire received utterance in a first step. The main part of  
the calculation is then focussed on the reanalysis of the utterance for  
extracting its meaning, intention and therefore for generating key-phrases and/  
or the keywords of the utterance. Keywords or key-phrases are parts or  
15 subunits of the utterances which carry the main importance of the message to  
be transported by the utterance. Consequently, the inventive method for  
recognizing speech saves calculational and estimation power by focussing on  
important parts of an utterance, namely the key-phrases and keywords, and on  
their generation, extraction and/or confidence estimation from the utterance.

20 For a dialogue system it is preferred that in the case of rejecting said utterance  
in its entirety a rejection signal is generated. In particular, a reprompting  
signal and/or an invitation to repeat or restart the last utterance is generated  
and/or output as said rejection signal. This is of particular advantage in a  
25 dialogue system as the user or current speaker is informed that his last  
utterance or speech phrase has not been recognized correctly by the  
recognizing system or method.

For performing the above mentioned rough estimate upon accepting and/or  
30 rejecting a received and/or recognized utterance a rough or simple confidence  
measure for the entire utterance is determined. This is of particular advantage  
in contrast to prior art methods for recognizing speech as these prior art  
methods generally calculate confidence measures which are based on each  
single word or subword unit within said utterance. Therefore, for the entire  
35 utterance prior art methods have to calculate and determine a relative large  
number of single word confidence measures.

- 1 Additionally, prior art methods for recognizing speech have then afterwards to perform an overall estimation to find a confidence for the whole utterance with respect to the set of single word confidence measures. In contrast to these prior art methods the inventive method calculates in the initial phase of recognition
- 5 a confidence measure for the whole utterance in its entirety and in a simple and rough manner. Only if on the basis of said whole utterance confidence measure an acceptance of the utterance and the recognized phrases thereof is suggested, further processing is initiated.
- 10 It is preferred to base said reanalysis on a sentence analysis, and in particular on grammar, syntax and/or semantic analysis or the like. These measures are useful as they are concentrated on extracting the intention and the meaning as well as on the extraction of the key-phrases or keywords of the utterance. In particular, in dialogue systems it is necessary that the method implemented in
- 15 the system is able to extract from the more or less complex received utterance the most important parts thereof so as to reduce the more or less complex utterance to its intention and meaning, in particular by collecting the key-phrases or keywords.
- 20 It is therefore of further advantage to form a relatively thorough estimation on whether the extracted key-phrases and/or keywords of the utterance can be accepted or have to be rejected in particular by the previous confidence measure.
- 25 In a particular advantageous embodiment of the inventive method for recognizing speech a detailed and/or robust confidence measure for each single key-phrase/keyword is determined for said thorough estimation of accepting/rejecting said key-phrases and/or keywords.
- 30 To further reduce the computational burden of the inventive method for recognizing speech the above described detailed and/or robust confidence measure for the derived key-phrases/keywords of the received and recognized utterance is only derived if within said step of deriving said key-phrase/keyword an indication and/or demand therefor is generated or does occur.
- 35
- Some of the basic ideas of the inventive methods for recognizing speech in contrast to prior art methods can be described and summarized as follows:



- 1 Confidence measures (CM) try to judge on how reliable an automatic speech  
recognition process is performed with respect to a given word or utterance. The  
confidence measure proposed in connection with the present invention is  
particularly designed for dialogue systems which have to deal with continuous  
5 speech input and which have to perform distinct actions based on data  
extracted and gathered from the input and recognized speech. The inventive  
method for recognizing speech combines various sources of information to  
judge if an input and recognized utterance and/or the particular selected  
words are recognized correctly.

10

- After a first step of recognizing the utterance in its entirety a simple, rough and  
very general confidence measure is computed and generated for the whole, i. e.  
entire utterance. If the recognized utterance is classified as being accepted the  
method turns to a further step of processing. Depending on the requirements of  
15 the method particularly implemented in a system a more detailed confidence  
judgement for the words or subword units which are of special importance can  
be generated on demand. These words or subword units of special importance  
are called key-phrases or keywords. The further processing steps, i. e. the  
reanalysis of the utterance, may explicitly ask for the calculation of the  
20 reliability of the key-phrases and/or keywords in the sense of a detailed and  
more robust confidence measure focussing on the corresponding single key-  
phrases or keywords.

- For the judgement of recognition quality in large vocabulary continuous speech  
25 dialogue systems a two-step system is therefore proposed. The first step of  
recognizing the utterance entirely and of calculating a simple confidence  
measure gives an indication if most of the utterance was recognized correctly.  
For such a classification, however, not every single word of the user input is  
equally important. The knowledge about the importance is usually not located  
30 within the information stored in the speech recognition system. It is therefore  
proposed to add an interface to the speech recognition subsystem that allows a  
following component to query specifically for the confidence of single words of  
the recognized utterance.

- 35 Therefore, after the analysis of the meaning or intention of the utterance in its  
entirety, an isolated word, more complicated and more robust confidence  
measure is applied to the isolated words or short phrases of special interest,

- 1 i. e. it is applied to the key-phrases or keywords of the utterance, in particular  
on demand of following speech recognition subsystems for entirely specifying  
the utterance.
- 5 If standard methods for the confidence measure judgement would be applied at  
this stage this would enlarge the computational burden. One could simply  
extend the approach developed so far for isolated words to continuous speech  
recognition and compute a very detailed confidence measure for each single  
word in the utterance. Since this would be very costly, the system response  
10 would be slowed down. For dialogue systems which have to respond fast to the  
input utterance of the user or speaker this is not acceptable. Therefore, the  
inventive method is proposed as follows.

The purpose of the first processing step of computing a rather simple  
15 confidence measure for the utterance is to aid the finding of the general struc-  
ture of the utterance. If this classification is done with high enough confidence,  
subsequent steps of proceeding can further process the received and recognized  
utterance. In these further processing steps the sentence or utterance is  
further analyzed so as to identify the important keywords of the sentence or  
20 utterance. On demand for these keywords a second more detailed and thorough  
confidence measure can be computed. Furtheron, additional and more sophisti-  
cated features that need a high amount of computational effort can be used in  
the second run to compute a confidence measure. Thereby, the expensive  
computational pathway is reduced and focussed to those locations of the  
25 utterance where it is really needed in the context of the application. This  
reduces the overall computational load and makes confidence estimation feasi-  
ble in small appliances.

For example, in a train time table information system the user utters "I want to  
30 go from Hamburg to Stuttgart". The intention of this utterance is to go from  
one city to another. For this information only the starting city and the  
destination have to be verified, whereas the rest of the sentence can be  
considered as filling phrases or "fillers". These filling phrases have not to be  
recognized with high accuracy as long as the intention of travelling from one  
35 point to another is known. Therefore, what is important is to verify the start  
city and destination. Therefore, according to the invention the computational  
load is focussed to these keywords, i. e. the start and destination of the

- 1 intended travel. Therefore, the second confidence measure is computed - if  
required - on start and destination only.

In other applications the speech recognizer outputs alternative word hypothe-  
5 ses arranged in a graph in order to cope with uncertainties and ambiguities.  
There exist many possible paths in the word graph each of which corresponds  
to a sentence hypothesis. The subsequent linguistic processor searches for the  
optimal path according to linguistic knowledge and to acoustic scores  
previously computed in the speech recognizer. During the search where the  
10 linguistic processor parallely explores several paths it may demand the  
confidence measure calculating module to score certain keywords. That means,  
at each following step a confidence measure can be queried. Which words are  
the keywords depends on the current stage of syntactic and semantic analysis  
within the underlying syntactic/semantic analysis.

15

The invention will be shown in more detail by means of a schematical drawing  
describing a preferred embodiment of the inventive method for recognizing  
speech.

- 20 Figure 1 describes by means of a schematical block diagram an embodiment of  
the inventive method for recognizing speech.

In a first step 11 continuous speech input is received as an utterance U and  
preprocessed. In step 12 a large vocabulary continuous speech recognizing  
25 process LVCSR is performed on the continuous speech input, i. e. the received  
utterance U or speech phrase so as to generate a recognition result in step 13.  
The recognition result of step 13 serves as an utterance hypothesis which is fed  
into step 14 for calculating a simple and rough confidence measure CMU for  
the entire utterance hypothesis of step 13. In the case of a rejection given by  
30 the confidence measure CMU of the whole utterance hypothesis a reprompt or  
invitation to repeat the utterance is initiated in step 20.

In the case of an acceptance of the utterance hypothesis a thorough sentence  
analysis is performed in step 15 so as to extract keywords in step 16. In a  
35 further step 17 it is calculated whether or not a confidence measure is  
necessary to evaluate the keywords. If a further evaluation on the reliability of  
the extracted keywords is necessary a thorough confidence measure CMK

- 1 calculation is demanded using time-alignment information called from the large  
vocabulary continuous speech recognizing unit of step 12. If no confidence  
measure CMK was necessary or the confidence measure CMK for the keywords  
was sufficient, the generated and extracted keywords and key-phrases are  
5 accepted. If the detailed confidence measure CMK was not sufficient the  
keywords are rejected and a reprompt is initiated branching the process to step  
20.

10

15

20

25

30

35

MÜLLER &amp; HOFFMANN

- 9 -

Sony International (Europe) GmbH

File: 52.744

16.11.2000

## Claims

EPO - Munich  
17

16. Nov. 2000

- 1 1. Method for recognizing speech,  
- wherein a received utterance (U) is subjected to a recognition process  
in its entirety,  
- wherein a rough estimation is made on whether or not said received  
5 utterance (U) is accepted or rejected in its entirety,  
- wherein in the case of accepting said utterance (U) it is thoroughly  
reanalyzed so as to extract its meaning and/or intention, and  
- wherein based on the reanalysis keywords and/or key-phrases are  
extracted from the utterance (U) essentially being representative for its  
10 meaning.
2. Method according to claim 1,  
wherein in the case of rejecting the utterance (U) a rejection signal is  
generated.
- 15 3. Method according to claim 2,  
wherein as said rejection signal a reprompting signal and/or in the case  
of a dialogue system an invitation to repeat/restart the last utterance (U) is  
generated and/or output.
- 20 4. Method according to anyone of the preceding claims,  
wherein for said rough estimation on accepting/rejecting the utterance a  
rough and/or simple confidence measure (CMU) for the entire utterance (U) is  
determined.
- 25 5. Method according to anyone of the preceding claims,  
wherein said reanalysis of the received utterance (U) is based on a  
sentence analysis, in particular based on a grammar, syntax, semantic  
analysis and/or the like.
- 30 6. Method according to anyone of the preceding claims,  
wherein a thorough estimation is made on whether or not said extracted  
keywords and/or key-phrases are accepted or rejected.

1     **7.**     Method according to claim 6,

          wherein for said thorough estimation on accepting/rejecting said key-phrases and/or keywords a detailed and/or robust confidence measure (CMK) for each single key-phrase or keyword is determined in particular on demand.

5

**8.**     Method according to claim 7,

          wherein a confidence measure (CMK) for the single key-phrase/keyword is determined only if in the step of deriving said key-phrase/keyword and indication therefore occurs so as to reduce the computational burden.

10

15

20

25

30

35

**Abstract****METHOD FOR RECOGNIZING SPEECH**

To increase the performance rate of large vocabulary continuous speech recognition applications it is suggested to first give only a rough estimation on whether or not recognized utterance (U) is accepted or rejected in its entirety. In the case of an acceptance of the utterance (U) a thorough reanalysis is performed afterwards to extract the meaning, intention, contained key-phrases/keywords and the confidence of the contained key-phrases/keywords. Therefore, the computational burden is focussed on the important sections of the utterance (U), namely on the key-phrases/keywords.

(Figure 1)

EPO - Munich  
17  
16. Nov. 2000

**THIS PAGE BLANK (USPTO)**



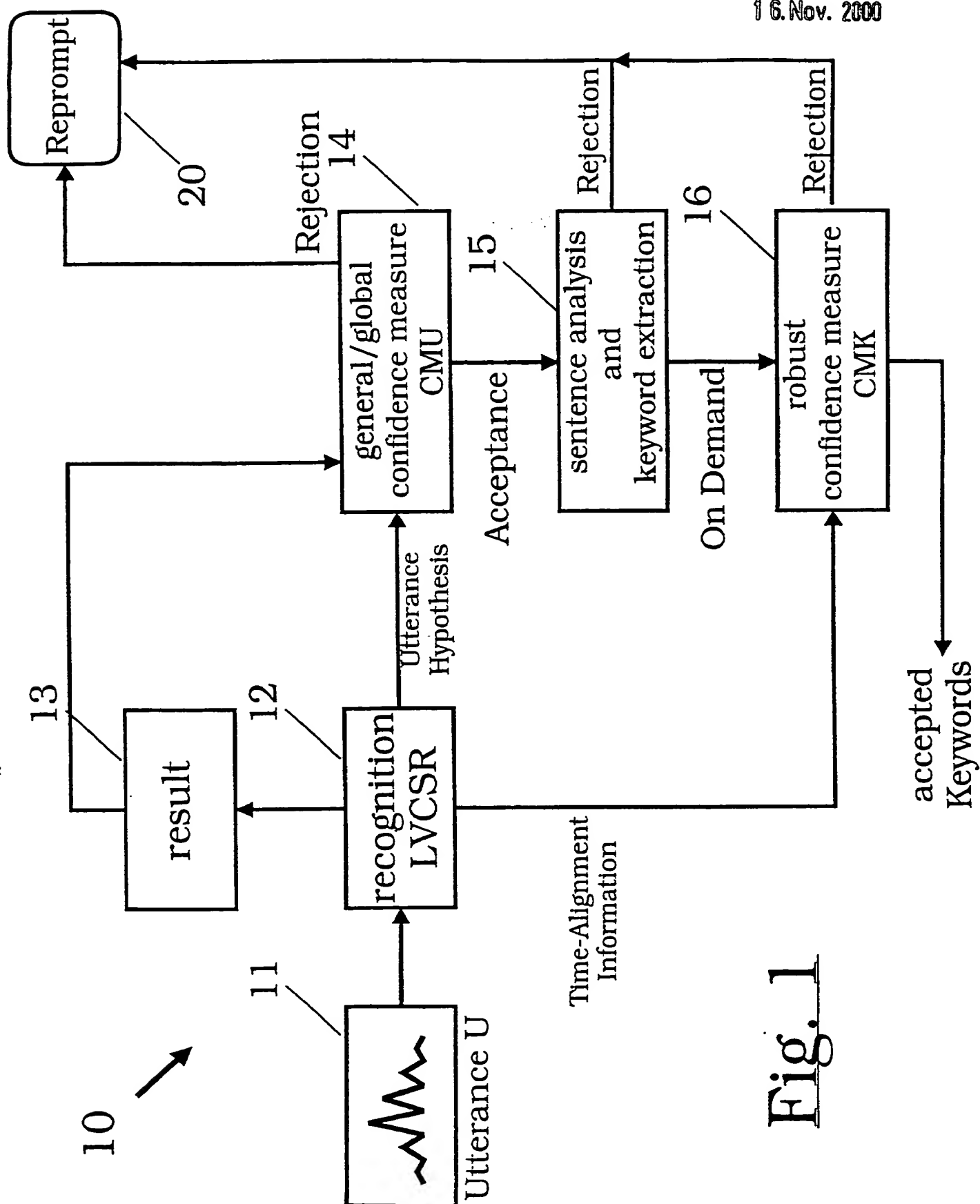


Fig. 1

**THIS PAGE BLANK (USPTO)**